



## Sensor Data Analysis

Prof. Dr. Dunja Mladenicić, Dr. Blaž Fortuna  
Jožef Stefan Institute and  
Jožef Stefan International Postgraduate School  
Slovenia


ailab.ijs.si 



## Information Age - Age of Analytics

*“This is the Information Age — everybody can be informed about **anything and everything**. There is no secret, therefore there is no sacredness. **Life is going to become an open book**. When your computer is more loyal, truthful, informed and excellent than you, you will be challenged. You do not have to compete with anybody. You have to compete with yourself.”*

[Y. Bhajan]

ailab.ijs.si 



## Overview

- Introduction
- Big Data from Sensors
- Algorithms for Big Data
- Techniques for (Big) Data Modelling
- References

ailab.ijs.si



## INTRODUCTION

ailab.ijs.si





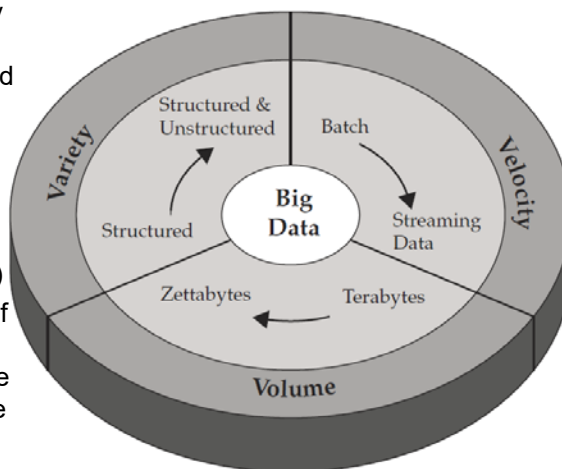
## Big Data

- 'Big data' is similar to 'Small data', but bigger
- ...but having data bigger requires different approaches – big spoon:
  - different techniques, tools, architectures
- ...with an aim to solve new problems
  - ...and old problems in a better way



## Characterization of Big Data: volume, velocity, variety (V3)

- **Volume** – data generated by machines, networks, social media, .... challenging to load and process (how to index, retrieve?)
- **Variety** – many sources and data types with different degree of structure (how to query semi-structured data?)
- **Velocity** – continuous flow of data requires real-time processing influenced by rate of data arrival (how to handle high rate?)



From "Understanding Big Data" by IBM

ailab.ijs.si



## The extended V3 of Big Data (Vn)

- 1. **Volume** (lots of data = “Tonnabytes”)
- 2. **Variety** (complexity, curse of dimensionality)
- 3. **Velocity** (rate of data and information flow)
- 4. **Veracity** (noise and outliers in data, need for verifying the inferred models)
- 5. **Variability** (variance in meaning)
- 6. **Venue** (location)
- 7. **Vocabulary** (semantics)
- 8. **Volatility** (how long is the data valid)

ailab.ijs.si



## Big Data in Data Science

Interdisciplinary field, combines methods from

- statistics, machine learning, analytics,
- visualization,
- reporting, business intelligence, expert systems,
- databases, data mining, big data

Process to transform hypotheses and data into actionable knowledge/predictions:

- Acquiring and managing the data
- Choosing the modeling techniques and writing the code
- Verifying the results

ailab.ijs.si





## Roles of People in Data Science

- Project sponsor
  - business interest, championing the project
- Client
  - domain expert, end user
- Data scientist
  - set and execute analytics, managing the project
- Data architect
  - data management and storing
- Operations
  - acquiring data, infrastructure management, deployment

Regular communication between data scientist with sponsor and with client, ensuring timely feedback

ailab.ijs.si



## Project Stages in Data Science

- Defining the goal
- Collecting and managing data
- Building the model
- Evaluating the model
- Results presentation
- Model deployment

Loop through the stages repeating as needed

ailab.ijs.si





## Big Data Transforming Business Models

- Volume of business data increasing (expected to double every 1.2 years)
- Different industries are using big data to transform business models and improve performance in many areas

Illustrative

Retail	Manufacturing
<ul style="list-style-type: none"> <li>Customer relationship management</li> <li>Store location and layout</li> <li>Fraud detection and prevention</li> <li>Supply chain optimization</li> <li>Dynamic pricing</li> </ul>	<ul style="list-style-type: none"> <li>Product research</li> <li>Engineering analytics</li> <li>Predictive maintenance</li> <li>Process and quality analysis</li> <li>Distribution optimization</li> </ul>
Financial services	Media and telecommunications
<ul style="list-style-type: none"> <li>Algorithmic trading</li> <li>Risk analysis</li> <li>Fraud detection</li> <li>Portfolio analysis</li> </ul>	<ul style="list-style-type: none"> <li>Network optimization</li> <li>Customer scoring</li> <li>Churn prevention</li> <li>Fraud prevention</li> </ul>
Advertising and public relations	Energy
<ul style="list-style-type: none"> <li>Demand signaling</li> <li>Targeted advertising</li> <li>Sentiment analysis</li> <li>Customer acquisition</li> </ul>	<ul style="list-style-type: none"> <li>Smart grid</li> <li>Exploration</li> <li>Operational modeling</li> <li>Power-line sensors</li> </ul>
Government	Healthcare and life sciences
<ul style="list-style-type: none"> <li>Market governance</li> <li>Weapon systems and counterterrorism</li> <li>Econometrics</li> <li>Health Informatics</li> </ul>	<ul style="list-style-type: none"> <li>Pharmacogenomics</li> <li>Bioinformatics</li> <li>Pharmaceutical research</li> <li>Clinical outcomes research</li> </ul>

Source: A.T. Kearney analysis

[Hagen et al., AT]

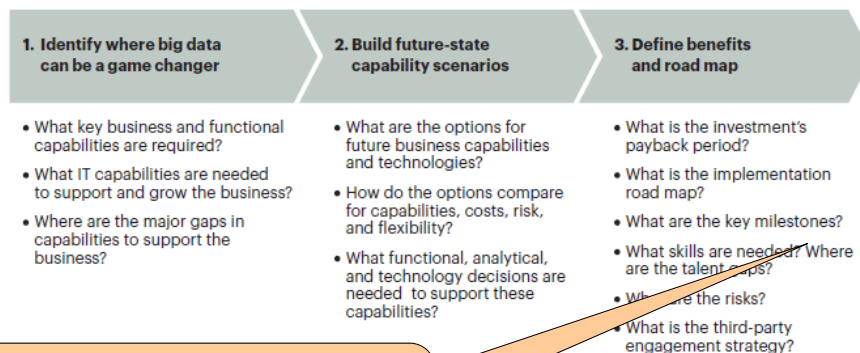
Big Data and the Creative Destruction of Today's Business Models 4



## Big Data Transforming Business Models (cont.)

### Getting started on the big data journey

Illustration



ETL – Extract data from various sources, Transform through workflow and data integration process, Load into database  
 ELT – Extract, Load the raw data and store potentially for longer time, Transform to use

More structured

SQL is structured query language. ETL is extract, transform, and load. ELT is extract, load, and transform.



A slide with a light blue and white background featuring a stylized butterfly logo in the top left corner. The text "When Big Data is really a hard problem?" is written in a bold, purple, sans-serif font in the center. Below the title, there are two main bullet points, each with a red dot, and several sub-bullet points, each with a yellow dot. In the bottom right corner, the text "ailab.ijs.si" is displayed next to a small black logo consisting of a grid of dots.

- ...when the operations on data are complex:
  - ...e.g. simple counting is not a complex problem
  - Modeling and reasoning with data of different kinds can get extremely complex
- Good news about big-data:
  - Often, because of vast amount of data, modeling techniques can get simpler (e.g. smart counting can replace complex model-based analytics)...
  - ...as long as we deal with the scale



## What matters when dealing with data?

- Data modalities, data operators, additional issues
- Research areas, such as IR, KDD, ML, NLP, Semantic Web, ... are sub-cubes within the data cube

ailab.ijs.si

## Truth or a random phenomenon?

Risk with “Big Data mining”

- we can “discover” patterns that occur by chance
- ...if you look in more places for interesting patterns than your amount of data will support, you are bound to find crap

**Bonferroni's principle**

we find a statistical pattern  
we are looking for

- a pattern is there by chance

“...truth is simple, straight and with a smile. You don't have to remember it. You have to say it. You know it and then you have to live it. It is so simple.” [Y.Bhajan]

ailab.ijs.si





## Meaningfulness of Analytic Answers

Calculate the expected number of occurrences of the pattern under the assumption that the data is random

### Illustrative example

- Find (unrelated) people who **at least twice have stayed at the same hotel on the same day** (can be different hotel each day)
  - $10^9$  people being tracked
  - 1000 days
  - each person stays in a hotel 1% of the time (1 day out of 100) – probability of staying in a hotel is 0.01
  - there are  $10^5$  hotels, capacity of a hotel is 100 people

If everyone behaves randomly (i.e., no conspiracy) will the data mining (by chance) detect anything suspicious?

Example taken from: Rajaraman, Ullman: Mining of Massive Datasets

ailab.ijs.si



## Calculation of patterns detected by chance

Event/pattern: 2 people on 2 days stay in the same hotel

- 2 people at the same day go to a hotel
  - a person stays in a hotel 1% of the time,  $0.01 * 0.01 = 10^{-4}$
- 2 people at the same day go to the same hotel ( $10^5$  hotels)
  - probability =  $10^{-4} * 10^{-5} = 10^{-9}$
- 2 people at the same day go to the same hotel, occurs twice
  - probability =  $10^{-9} * 10^{-9} = 10^{-18}$

$$\binom{n}{2} = \frac{n * (n - 1)}{2} = \frac{n^2 - n}{2} \approx \frac{n^2}{2}$$

### Random behavior

- Choose 2 people from  $10^9$  and choose 2 days from  $10^3$ 
  - ways to choose:  $10^{18}/2 * 10^6/2 = 5 * 10^{17} * 5 * 10^5 = 25 * 10^{22}$
- Event probability - expected number of “suspicious” pairs of people in random data (out of  $10^9$  people) = 250 000 (!)
  - $25 * 10^{22} * 10^{-18} = 25 * 10^4 = 250\,000$
  - ... too many combinations to check – we need to have some additional evidence to find “suspicious” pairs of people in some more efficient way

Example taken from: Rajaraman, Ullman: Mining of Massive Datasets

ailab.ijs.si





## Variation: $10^7$ people being tracked instead of $10^9$

### Random behavior

- Choose 2 people from  $10^7$  and choose 2 days from  $10^3$ 
  - ways to choose:  $10^{14}/2 * 10^6/2 = 5 * 10^{13} * 5 * 10^5 = 25 * 10^{18}$
- Event probability - expected number of “suspicious” pairs of people in random data of  $10^7$  people
  - $25 * 10^{18} * 10^{-18} = 25$

Example taken from: Rajaraman, Ullman: Mining of Massive Datasets

ailab.ijs.si



## Big Data from Data Stream

### Data stream is a common source of **big data**

- web logs, social media, stock market, sensor networks,...
- **Data stream management**
  - Problematic are blocking query operators – need the entire input to produce any result (eg, sort, sum, max)
  - use approximations, sampling, window of data
- **Data stream processing**
  - Maintain simple statistics on stream (mean, standard deviation)
  - Use time window:
    - sliding (fixed size – eg. the last 100 values),
    - landmark (fixed start – eg. from the start of the day)
    - tilted (recent data in more details – eg, last hour in 15 mins, last day in 24 hours, last month in days, last year in months)

ailab.ijs.si





## Sensor networks

- Networks of small sensing devices distributed over locations
- Capability to sense, process, act, communicate
- Sensor
  - Equipped with memory, processing capability, communication with neighbors,
  - Constrained on resources (energy, memory, computational speed, bandwidth)
  - Used for monitoring (eg., traffic), tracking (objects), controlling (production)
  - Producing big data: continuous flow of sensor readings often at high speed, in dynamic and time-changing environment, large number of sensors on different locations

ailab.ijs.si



## Sensor Data vs. Traditional Data

- Sample of the population – data from continuous stream
- Noisy – sensing equipment requires data cleaning
- Data duplication (similar environmental conditions over large area with many sensors)
- Spatial and temporal attributes play a major role
- Data processing on network nodes (energy & bandwidth limit transferring of all the data to a central site)
  - Distributed processing of queries, real-time data cleaning, energy efficiency
  - Limited computational resources on sensing nodes (limited bandwidth, processing power and memory, low-power batteries, scaling to many sensors, robust to noise)
  - Multi-level data modeling - local models differ from a global model

ailab.ijs.si





## Querying sensor data

- Sensor node – mini data repository
- Sensor network – database distributed across sensor nodes
- Query is sent to sensor network
  - through gateway – a special purpose node
  - forwarded hop-by-hop from the gateway to sensor nodes - query broadcasting (can be selective, eg., only to the nodes at requested location or only nodes that measure temperature)
- Sensor nodes probe their sensing devices and propagate sensor readings back to the gateway
- In-network processing on nodes – uses less energy to send aggregates/approximations than wireless data transmission

ailab.ijs.si

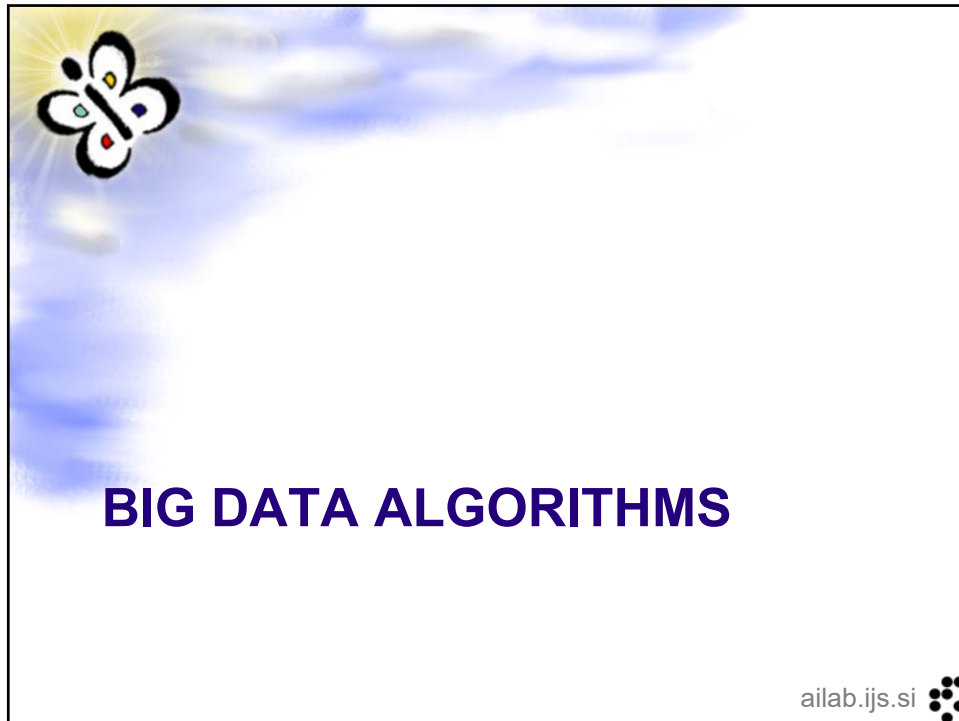


## Queries in sensor networks

- One-shot vs. long-running queries
  - eg, location of empty parking spot vs. hourly monitoring of temperature
- All data vs. aggregate queries (eg, daily average)
- Time-based vs. event-based queries
  - eg, on temp >40 start sending temp every second
- Accurate vs. approximate queries
- Urgent vs. delay-tolerant
  - eg, intruder detection vs. average hourly temperature
- Pull vs. push
  - explicit query vs. sensor nodes sending data regularly

ailab.ijs.si





A slide titled "Types of tools typically used in Big Data scenarios" in bold, purple capital letters. In the top left corner, there is a stylized butterfly logo with colorful spots. In the bottom right corner, the text "ailab.ijs.si" is displayed next to a small cluster of black dots. The background of the slide features a soft, blue and white cloud-like pattern.

- Where the processing is **hosted**?
  - Distributed Servers / Cloud (e.g. Amazon EC2)
- Where the data is **stored**?
  - Distributed Storage (e.g. Amazon S3)
- What is the **programming model**?
  - Distributed Processing (e.g. MapReduce) - very simple operations on large volume of data
- How the data is **stored & indexed**?
  - High-performance schema-free databases (e.g. MongoDB)
- What **operations** are performed on data?
  - Analytic / Semantic Processing (e.g. R, OWLIM)



## Big Data Analytics

- **Smart sampling** of data
  - ...reducing the original data while not losing the statistical properties of data
- **Finding similar items**
  - ...efficient multidimensional indexing
- **Incremental updating** of the models
  - (vs. building models from scratch)
  - ...crucial for streaming data
- **Distributed linear algebra**
  - ...dealing with large sparse matrices

ailab.ijs.si



## Analytical operators on Big Data

- On the top of the previous ops we perform usual data mining/machine learning/statistics operators:
  - **Supervised** learning (classification, regression, ...)
  - **Semi-supervised** learning (co-training, active learning,...)
  - **Unsupervised** learning (clustering, different types of decompositions, ...)
- ...we are just more careful which algorithms we choose (typically linear or sub-linear versions)

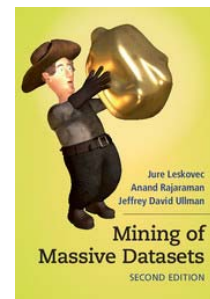
ailab.ijs.si





## ...guide to Big-Data algorithms

- An excellent overview of the “Big Data” algorithms is the book “**Leskovec, Rajaraman, Ullman: Mining of Massive Datasets**”
  - Downloadable from: <http://www.mmids.org/>
  - Associated MOOC (from Oct 2014): <https://www.coursera.org/course/mmids>



## “Big Data Research” Journal



- Elsevier started new “Big Data Research” journal
  - <http://www.journals.elsevier.com/big-data-research/>
- Special issues
  - [Visions on Big Data](#)
  - [Special Issue on Computation, Business, and Health Science](#)
  - [Big Data, Analytics, and High-Performance Computing](#)
- Recent Articles
  - [Finding the Best Classification Threshold in Imbalanced Classification](#)
  - [Analysis of a Network IO Bottleneck in Big Data Environments Based on Docker Containers](#)
  - [Practical Identification of Dynamic Precedence Criteria to Produce Critical Results from Big Data Streams](#)







## Sampling on Big-Data

### Sampling

- Deals with velocity and volume
- Enables off-line data analysis
- Enables performing expensive operations (eg, join of two streams via join of two samples)

**Reservoir sampling** – maintaining a sample of fixed size by probabilistically replacing an old element by a new one

**Sampling from different data sources** – depending on the kind of queries to be asked decide whether to consider info. about the data source

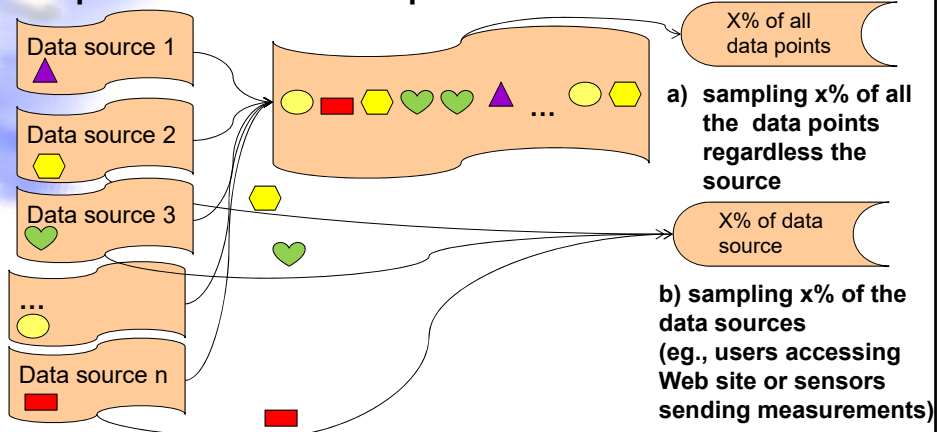
**Sampling  $a/b$  fraction of the data** – hash data into  $b$  buckets and decide whether to store the data point based on the calculated value of the hash function

ailab.ijs.si



## Sampling different data sources

Depends on the kind of queries that will be asked




Example: average number of Web pages revisited by the same user?

- requires b) sampling - all the data for x% of the users

ailab.ijs.si





## Sampling data using hashing

Diagram illustrating the process of sampling data using hashing:

**Data source(s)** → **hashing** → **buckets**


The buckets are numbered 1, 2, ..., a, a+1, ..., b. Each bucket contains a specific data element (e.g., a heart, a red square, a yellow hexagon, a purple triangle, a crescent moon, and a star).

- use a hash function to hash the key components of the data stream (eg., username),
- based on the value of the function decide whether to store the current data or not

**Example:**

- hash the username to  $b$  buckets,
- if the user falls into one of the first  $a$  buckets store the data

ailab.ijs.si



## Finding similar items – set intersection

Approach as a problem of finding sets with large intersections - Jaccard similarity:

$$\text{set\_intersection} / \text{set\_union}$$

- a) estimate by making
  - random permutations of elements in  $\text{set\_intersection}$  and  $\text{set\_union}$
  - compare the first elements (they are equal with probability of Jaccard similarity)
- b) instead of random permutations use hashing

Focus on similarity between the promising pairs of items

- eg., usernames with the same hash value, documents of the same length

ailab.ijs.si



## Finding similar items – across similarity

Approach by estimating similarity measure instead of calculating the exact value

- estimate cosine similarity by estimating an inner product of two vectors
  - multiply each vector with some random vector  $N(0,1)$
  - if both results are positive or both are negative assume the original vectors are similar

Example problem

- Similarity of documents (plagiarism, mirror Web pages, news articles from the same source)
- Collaborative filtering for movie/book/... recommendation

ailab.ijs.si



## Storing Big Data

- Data arriving in streams, rapidly so it is not feasible to store all the data
  - Eg., measurements of sensors at different locations – even if one stream is not of high speed, there is multitude of streams
- **What to store** depends on the queries that will be asked
  - **Standing query** (event pattern)
    - trigger an alarm, perform an operation on each arrival of a data point (eg., average the last 100 readings of sensor), report max. temperature so far
  - **Ad-hoc query**
    - Store **sliding window** of the last  $n$  data points
      - eg., the last 10 values of wind speed
    - Store **the last  $t$  time units** readings
      - eg., wind speed during the last hour,
      - eg., the number of unique users on the Web site in the past month
      - store the complete stream for the last month with the time stamp, remove the old data as new arrives

ailab.ijs.si





## Storing Data on more Machines

- Most big data have inherent temporal and/or spatial dimension
  - Data with time dimension should be stored and processed at least in a partial temporal ordering
  - Distributed storing of the data should consider the kind of queries that will be asked
    - if we want different type of queries i.e. over time and over location the data can be replicated to improve efficiency (and provide redundancy over potential hardware failure)
- A cluster of 10 machines is 10 times more likely to require a service than one machine

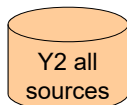
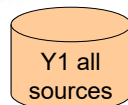
ailab.ijs.si



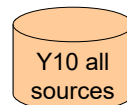
## Storing data on more Machines - Example

Example: 10 years of observations collected at 15s intervals from 1000 sensor sites can be stored on 10 machines:

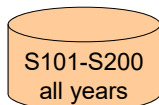
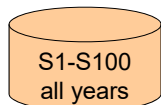
- a) All observations for each year on one machine (eg., to return average value for the last year of all sensors)
- b) All observations for 100 sensors on one machine (eg., to make analysis for one sensor over 10 years)



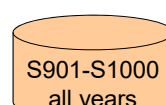
...



a) Year –centric




...




b) Source-centric

ailab.ijs.si






## Data Analytics Publicity - “Slaves of the algorithms”



The Economist, May/June 2013

Data analysis to **predict financial outcome of a movie** give the main actors – suggests paying some actors not to be on the movie helps avoiding making financial loss ;-)

ailab.ijs.si 



## TECHNIQUES FOR BIG DATA MODELING

ailab.ijs.si 



## Statistical Inference for Big Data

- Assess quality of statistical estimators via
  - Bootstrapping - resampling with resamples of size comparable to the original dataset - infeasible with massive data
  - Subsampling – resample size  $<$  the original dataset
- Bag of little bootstraps (BLB)
  - combining bootstrap and subsampling to obtain a robust, computationally efficient means of assessing estimator quality
  - Idea:
    - Average results of bootstrapping multiple small subsets (limit to  $b$  distinct points,  $b < n$ ) instead of working with large sets with  $0.632n$  distinct points in bootstrapping
    - Storing only  $b$  counts

[M. Jordan, 2012]

ailab.ijs.si



## Knowledge Acquisition

### Knowledge Valut from Google

- use supervised machine learning to automatically construct a Web-scale probabilistic knowledge base (1.6B RDF triplets subj-pred-obj with confidence score, 271M confidence  $>0.9$  of them 1/3 are new)
- Combine knowledge:
  - extraction from Web content (analysis of text, tabular data, page structure, human annotations) and
  - prior knowledge derived from existing knowledge repositories
    - uses Freebase but can be some other large-scale knowledge base (Wikipedia, DBPedia, NELL, YAGO, Microsoft's Satori, Google's Knowledge Graph)
- Calculate probabilities of fact correctness based on
  - other facts from knowledge base (using link prediction) and
  - extractors confidence (#sources and #extractors supporting the fact)

[Dong et al, KDD-2014]

*The acquisition of knowledge is always of use to the intellect, because it may thus drive out useless things and retain the good. For **nothing can be loved or hated unless it is first known.*** [Leonardo da Vinci]

ailab.ijs.si





## Cleaning Data

Empirical Glitch Explanations - concise, multi-dimensional descriptions of subsets of potentially dirty data

- Integrating large volumes of data from different sources brings inconsistency
- Data quality constraints to remove inconsistent data can be too strict
- Identify legitimate data and refine data quality constraints

Findings: significant portions of data that seem to violate constraints but have valid explanations and can be released back into the clean pool of data

[Dasu et al., KDD-2014]

ailab.ijs.si



## Clustering on streams

- BFR algorithm – k-means variant **assuming clusters are normally distributed** around the centroid
  - Instead of storing points, store summaries of the clusters + summaries of isolated mini clusters + outliers
- CURE – instead of centroid using a collection of **representative points**
  - Cluster a small sample of data to choose representative points, move representative points towards centroids, merge clusters with similar representatives
  - Assign all other points to one of the clusters based on similarity to representatives
- Clustering on a **sliding window** – assumes we are interested in clustering of the last  $m$  points

[Rajaraman et al, 2014, 255-266]

ailab.ijs.si





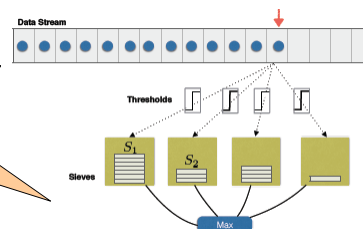


## Data Summarization on the Fly

- Select a subset of  $k$  most representative data points from a stream
  - cardinality-constrained submodular maximization
    - adding  $e$  in context of  $A$  helps at least as much as adding  $e$  in context of superset  $B$  of  $A$
- Keep the best  $k$  points in memory, as a new point arrives, if better than an existing point do replace
  - $m = \max \text{ quality of single } e, m \leq \text{opt} \leq k * m$

Constant factor approximation guarantees, no assumptions on the data stream, requires only a single pass, only  $O(k \log k)$  memory and only  $O(\log k)$  update time, assuming nothing but monotone submodularity

Algorithm



[Badanidiyuru et al., KDD-2014]



## Online Learning for Distributed Mining

- **Heterogeneous data from distributed sources**
- **Linear regression problems on feature-distributed data**

Exploiting correlations between local learners to reduce info. exchange and computational complexity

- **Ensemble learning (linear regression) with multiple local online learners – each has limited data access (feature-based partitioning)**
- **Local learners grouped based on correlation of their models**
- **Based on feedback from the ensemble learner**
  - cooperative updating of the correlated models,
  - independent updating of uncorrelated models

[Zhang et al., SIGMETRIC-2014]

ailab.ijs.si





## Querying Big Data

- Adaptive indexing – build parts of the index as needed for the users pose queries [Zoumpatianos et al, SIGMOD-2014]
- Querying Big Data – algebraic layer of complex query processing
  - Similarity of complex objects, combining semi-structured and unstructured search [Novikov et al., CompSysTech-2012]
- SQL RDBMS for Big data
  - Smaller transaction volume with large no. of rows per operation → idea: data stored column wise and compressed (vs. traditional RDBMS with row based storage, uncompressed)
- Incrementally leaf classifier Error Adaptive
  - Balancing data [Yang]

SYMBOL	DATE	TIME	PRICE	VOLUME	ETC
HPQ	05/13/11	01:02:03 PM	40.07	100	...
IBM	05/13/11	01:02:03 PM	171.22	10	...
AAPL	05/13/11	01:02:03 PM	338.02	9	...
GOOG	05/13/11	01:02:04 PM	524.03	150	...
HPQ	05/13/11	01:02:05 PM	30.07	40	...
AAPL	05/13/11	01:02:07 PM	338.02	20	...
GOOG	05/13/11	01:02:07 PM	524.02	40	...

SYMBOL (8K Distinct)	DATE (250yr)	VOLUME
GOOG (x18M)	05/13/2011 (x150K)	...
		20
		150
		40
HPQ (x22M)	05/13/2011 (x220K)	...
		20
		100
		40
IBM (x19M)	05/13/2011 (x150K)	...
		200
		10
		15



## The Era of Big Data

- In science available massive streams of data
    - astronomy, high-energy physics, ecology, genetics and molecular biology
  - In technology, personalization
    - data on fine-grained aspects of human behavior permitting the development of new services that are tailored to individuals
- Big Data requires consideration of
- systems issues
    - how to store, index and transport data at massive scales; how to exploit parallel and distributed platforms,
  - statistical issues
    - how to cope with errors and biases of all kinds; how to develop models and procedures that work big data,
  - algorithmic issues
    - how to perform computations using resources that scale as linear or sub-linear functions
  - legal, commercial and social issues

[M. Jordan, 2011]

ailab.ijs.si





## Big Data for Business

Be smart when using Big Data, combine different activity of the mind to achieve efficient utilization of:

- data and input (analytical mind)
- people and time (administrative mind)
- funds (financial mind)
- taking all into account in making executive decisions (executive mind)

[Sadhana Singh, 2015]

Data is a valuable asset in business, but before going for using (big) data (executive), check:

- What is the business problem or goal?
- Is the available data suitable? (analytical)
- What is the expected return on investment? (financial)
- Can we do it with the available resources timely?

ailab.ijs.si



## Big Data for Business (cont.)

- Volume, Velocity, Variety of Big Data requires tradeoff on data freshness, query response time, data quality and answer quality
- Research challenges:
  - Support feature engineering and selection (eg., scoring individual features)
  - Learning from partially labeled data (eg., active learning)
  - Managing missing data across heterogeneous stream
  - Combining offline and online learning
  - Interactive and collaborative mining
  - Visualization of Big Data
  - Privacy and transparency
- Approach Big Data in scientific, practical and economic fashion

[Gopalkrishnan et al., BigMine-2012]

ailab.ijs.si







## Open Source Big Data Tools

### Infrastructure:

- Kafka [<http://kafka.apache.org/>]
  - A high-throughput distributed messaging system
- Hadoop [<http://hadoop.apache.org/>]
  - Open-source map-reduce implementation
- Storm [<http://storm-project.net/>]
  - Real-time distributed computation system
- Cassandra [<http://cassandra.apache.org/>]
  - Hybrid between Key-Value and Row-Oriented DB
  - Distributed, decentralized, no single point of failure
  - Optimized for fast writes

ailab.ijs.si



## Open Source Big Data Tools Machine Learning

- Mahout
  - Machine learning library working on top of Hadoop
  - <http://mahout.apache.org/>
- MOA
  - Mining data streams with concept drift
  - Integrated with Weka
  - <http://moa.cms.waikato.ac.nz/>

### Mahout currently has:

- Collaborative Filtering
- User and Item based recommenders
- K-Means, Fuzzy K-Means clustering
- Mean Shift clustering
- Dirichlet process clustering
- Latent Dirichlet Allocation
- Singular value decomposition
- Parallel Frequent Pattern mining
- Complementary Naive Bayes classifier
- Random forest decision tree based classifier

ailab.ijs.si





## ...about anything and everything

- Big Data is everywhere, we are just not used to deal with it
- The “Big Data” hype is very recent
  - ...growth seems to be going up
  - ...evident lack of experts to build Big Data apps
- Can we do “Big Data” without big investment?
  - ...yes – many open source tools, computing machinery is cheap (to buy or to rent)
  - ...the key is knowledge on how to deal with data
  - ...data is either free (e.g. Wikipedia) or to buy (e.g. twitter)

ailab.ijs.si



<http://ailab.ijs.si>

**Artificial Intelligence Laboratory**

**Menu**

- Home
- Contact data
- Overview
- Projects
- People
- Organized events
- Public media
- Publications
- Internal seminar
- News
- Tools

**Meta**

Log in

**JSI Brown Bag Seminar**

**SOLOMON SEMINARJI**

**videolectures.net**  
exchange ideas & share knowledge

**Home**

English

Jožef Stefan Institute

**Artificial Intelligence Laboratory**

Jožef Stefan Institute

Ljubljana Slovenia

Promotional video

The Artificial Intelligence Laboratory is concerned mainly with research and development in information technologies with an emphasis on artificial intelligence. The main research areas are the following: (a) data analysis with an emphasis on text, web and cross-modal data, (b) scalable real-time data analysis, (c) visualization of complex data, (d) semantic technologies, (e) language technologies.

In collaboration with the Department of Communication Systems (E6) and Centre for Knowledge Transfer in Information Technologies (CTI) we have established a Cross-department laboratory for wireless sensor networks (SensorLab). The goal is to combine technologies for (a) sensor data acquisition, (b) communication between sensor devices, (c) statistical real-time data analysis, (d) semantic technologies, and to enable a wide range of research and development in different application areas, such as energy, ecology, transport, security, and logistics.

The Artificial Intelligence Laboratory puts special emphasis on the promotion of science. In collaboration with the Centre for Knowledge Transfer in Information Technologies (CTI) we are developing the VideoLectures.NET educational portal and organizing the national ACM competition in Computer Science (in Slovene).

The Artificial Intelligence Laboratory has a well-established collaboration with a number of academic and commercial organizations, some members of the Laboratory are involved with Stanford University, University College London, Jožef Stefan International Postgraduate School and companies Quintelligence, Cycorp Europe, LifefitLive, Madro Olo and Envisage.

**news**

Computer is human on memory, creativity, emotions, consciousness (artificial intelligence Slovene).

**Cross-lingual News**

iOS application targets professionals and the public.

**Cross-lingual Global Monitoring**

Dunja Mladenic at Ca Mation University (v)

**Media Monitoring**

Dunja Mladenic at Na on Event Registry.

**The Zois Ambassador**

awarded to Prof. Dr.



## Requirements for this class

- Attendance of the lectures and independent work on the assigned seminar following the provided instructions
- Report on the results of the project work to be sent via e-mail by 15.02.2017 to Blaz.Fortuna@ijs.si
  - 5-10 pages report
- Presentation of the seminar work on 2.03.2017 11:00
  - 5-10 slides presentation (10-15 minutes presentation)
- Oral exam on 2.03.2017

Notice for the next class 14.12.2016 11:00-15:00

- please bring your laptop
- check <https://github.com/blazf/mpsPractice> in advance for details on what software you need installed